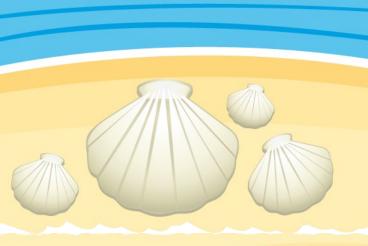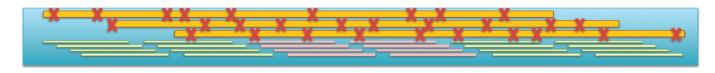# PacBio Assembly Softwares

End to End

**SMRT Hybrid**: Refers to the hybrid *de novo* assembly of error-corrected PacBio Continuous Long Read ("CLR") data (lower accuracy) with a second data type with higher accuracy - either PacBio Circular Consensus Sequence reads ("CCS") data or 2$^{nd}$ generation short-read data.

**SMRT Scaffolding**: Refers to using PacBio CLR to scaffold existing contigs generated from short-read data.

**SMRT *de novo***: Refers to the assembly of PacBio CLR data **only.**
**HGAP**: Using all PacBio data only.

**SMRT Gap Filling**: Refers to using PacBio CLR to fill gaps in existing mate pair-based scaffolds

## PacBio-only

| | | |
|---|---|---|
| ♥ | HGAP | A workflow to first preassemble reads, assemble the preassembled reads using Celera® Assembler, then polish using Quiver.<br>• Supports up to 100 Mb from SMRT Portal, which is part of SMRT Analysis.<br>• Larger genomes are possible from the command line using either smrtpipe.py or the Makefile-based smrtmake. |
| ♥ | Falcon | An experimental diploid assembler, tested on multi Gb genomes. 2014 AGBT presentation by Jason Chin. |
| | Canu | A fork of the Celera Assembler designed for high-noise single-molecule sequencing. |
| | Celera® Assembler | Celera® Assembler 8.1 now offers a way to directly assemble subreads. |
| | Sprai | A preassembly-based assembler that aims to generate longer contigs. |

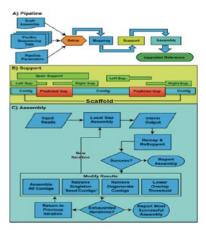| Hybrid | | |
|---|---|---|
| | pacBioToCA | An error correction module in Celera® Assembler originally designed to align short reads to PacBio reads and generate consensus sequences. These error corrected reads can then be assembled by Celera® Assembler. |
| ♥ | ECTools | A set of tools that uses contigs instead of short reads for correction. |
| | SPAdes | A short read assembler that added PacBio hybrid assembly support as of version 3.0. |
| | Cerulean | Cerulean starts with an assembly graph from ABySS and extends contigs by resolving bubbles in the graph using PacBio long reads. Was successfully run on genomes <100 Mb. |
| ♥ | dbg2olc | dbg2olc uses Illumina contigs as anchors to build an overlap graph with PacBio reads, allowing very fast performance. |

## Gap Filling

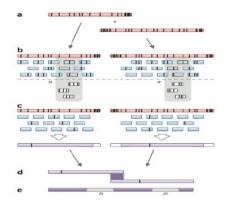| | | |
|---|---|---|
| ♥ | PBJelly 2 | PBJelly upgrades genomes by using PacBio reads to fill in gaps in scaffolds. Has been shown to work with genomes >1 Gb. Part of the PBSuite of applications including PB Honey. See also PAG 2014: Kim Worley, "Improving Genomes using Long Reads and PB Jelly 2 |

# PBJelly



**Gap Filling
and Assembly Upgrade**

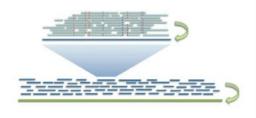English *et al* (2012)
*PLOS One.* 7(11): e47768

# PacBioToCA
# & ECTools



**Hybrid/PB-only Error
Correction**

Koren, Schatz, *et al* (2012)
*Nature Biotechnology.* 30:693–700

# HGAP & Quiver



$$\Pr(\mathbf{R} \mid T)$$
$$\Pr(\mathbf{R} \mid T) = \prod_{k} \Pr(R_k \mid T)$$

| Quiver Performance Results *Comparison to Reference Genome (M. ruber ; 3.1 MB ; SMRT® Cells)* | | |
|---|---|---|
| | Initial Assembly | Quiver Consensus |
| QV | 43.4 | 54.5 |
| Accuracy | 99.99540% | 99.99964% |
| Differences | 141 | 11 |

**PB-only Correction &
Polishing**

Chin *et al* (2013)
*Nature Methods.* 10:563–569

**< 5x**     **PacBio Coverage**     **> 50x**

# SMRT Portal Overview

# SMRT VIEW



Modification density by mod, variants, and separate tracks per motif

IPD ratio by base position and strand

# 下机数据文件夹内容

```
.
├── Analysis_Results
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.1.bax.h5  ♥
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.1.bax.h5.debug
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.1.log
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.2.bax.h5  ♥
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.2.bax.h5.debug
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.2.log
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.3.bax.h5  ♥
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.3.bax.h5.debug
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.3.log
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.bas.h5  ♥
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.bas.h5.debug
│   ├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.sts.csv
│   └── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.sts.xml
├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.1.xfer.xml
├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.2.xfer.xml
├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.3.xfer.xml
├── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.mcd.h5
└── m140913_102942_42208_c100658462550000001823128911271480_s1_p0.metadata.xml  ♥
```

# SGE

集群，简单的说即设置
一台电脑用作主控主机
（ qmaster ），收集集
群信息、分配任务、调
节负载均衡等；设置多
台电脑用作执行主机

qsub – 此命令是将作业提交到 Sun Grid Engine 系统的用户界面。
qhost – 此命令显示 Sun Grid Engine 执行主机的状态信息。
qdel
qstat
...

# SMRT Analysis 的安裝

```
Usage: smrtanalysis_2.3.0.140936.run [--help] [--rootdir dir] \
        [--skip-extract|--no-extract] [--extract-only] \
        [--update] \
        [-p|--patchfile patchfile] \
        [-a|--addonfile addonile] \
        [--start-origcmd args... --end-origcmd] \
        [--otherargs] \
        [-- thisprog_args... [-- subprog1_args... [subprog2_args...]]]
    --rootdir      -- smrtanalysis root directory (default: ./smrtanalysis)
    --skip-extract  -- skip the tarball extraction (use previous extract)
    --no-extract    -- same as --skip-extract
    --extract-only  -- only extract the tarball, do not invoke
                    installer or upgrader
    --update        -- update from an existing install (internal option
                    only)
    --patchfile     -- patch file to apply during install/upgrade
    --addonfile     -- addon file to apply during install/upgrade
    --start-origcmd -- original args to parent program, until
                    until --end-origcmd arg
    --otherargs     -- unrecognized args passed to subprogs until
    -- args...      -- force args to be handled by this program
    -- -- args...   -- force args to be handled by this subprog1
                    (first level subprog)
    -- -- -- args... -- force args to be handled by this subprog2
                    (second level subprog)
                    recognized
    --help          -- print this usage
    -- --help       -- print this usage
    --helpall       -- print this usage and usage of subprogs
```

```
[root@node2 opt]# pwd
/opt
[root@node2 opt]# ls
FALCON-integrate    sge             smrtanalysis_2.3.0.140936.run
PBSuite_14.1.15     smrtanalysis    test
[root@node2 opt]# ./smrtanalysis_2.3.0.140936.run --help
```

```
[root@node2 opt]# ./smrtanalysis_2.3.0.140936.run --help
Usage: smrtanalysis_2.3.0.140936.run [--help] [--rootdir dir] \
            [--skip-extract|--no-extract] [--extract-only] \
            [--update] \
            [-p|--patchfile patchfile] \
            [-a|--addonfile addonile] \
            [--start-origcmd args... --end-origcmd] \
            [--otherargs] \
            [-- thisprog_args... [-- subprog1_args... [subprog2_args...]
]]
        --rootdir        -- smrtanalysis root directory (default: ./smrtana
lysis)
        --skip-extract   -- skip the tarball extraction (use previous extr
act)
        --no-extract     -- same as --skip-extract
        --extract-only   -- only extract the tarball, do not invoke
                            installer or upgrader
        --update         -- update from an existing install (internal opti
on
                            only)
        --patchfile      -- patch file to apply during install/upgrade
        --addonfile      -- addon file to apply during install/upgrade
```

```
[root@node2 bin]# smrtpipe.py
-bash: smrtpipe.py: command not found
[root@node2 bin]# pwd
/opt/test/install/smrtanalysis_2.3.0.140936/smrtcmds/bin
[root@node2 bin]# ls
java  mono  perl  python  smrtpipe  smrtshell  smrtwrap
[root@node2 bin]# ./smrtshell
(smrtshell-2.3.0) [root@node2 bin]# smrtpipe.py
Usage: smrtpipe.py [--help] [options] dataUrl &> smrtpipe.err

smrtpipe.py: error: Expected 1 argument
(smrtshell-2.3.0) [root@node2 bin]#
```

# SMRT Portal Overview

# 数据导入



Import and Manage



**Manage Protocols**

Create and edit standard protocols for secondary analysis jobs in SMRT Portal.

**Manage Reference Sequences**

Import and manage reference sequences for resequencing and visualization with SMRT View.

**Import SMRT Cells**

Import raw data from SMRT cells for analysis in SMRT Portal.

**Import SMRT Pipe Jobs**
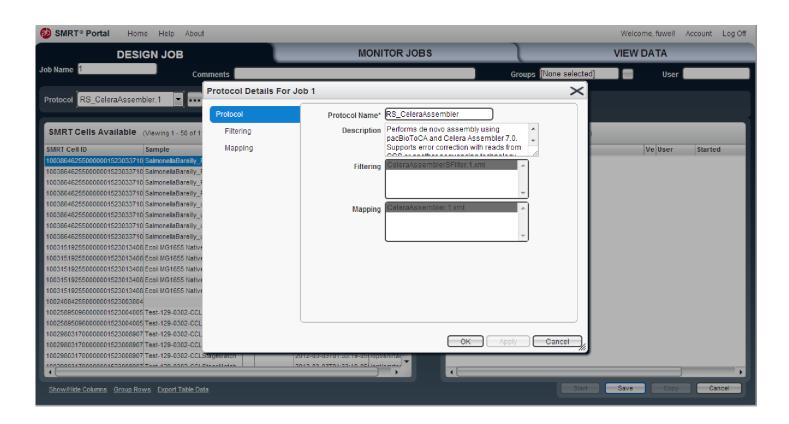
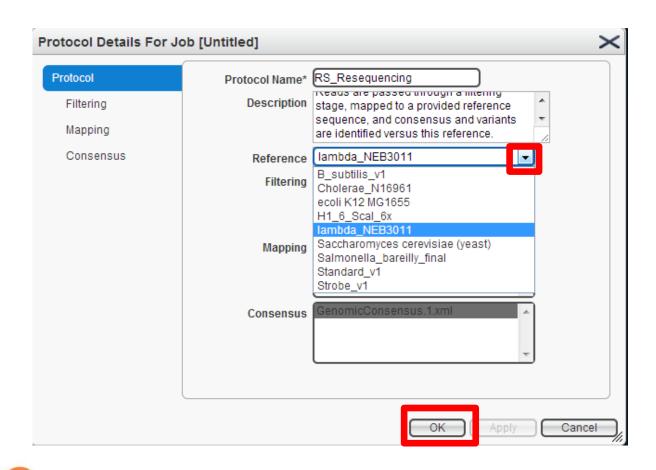Import SMRT Pipe jobs for display in SMRT Portal.

# 新建任务

# Assembly

# Resequencing

# Ref selection

# Parameter

**Protocol Details For Job [Untitled]** ✕

Protocol

**Filtering**

Mapping

Consensus

**SFilter v1**

**Minimum Readlength** 50

**Minimum Subreadlength** 50

**Minimum Read Quality** 0.75

**Description:** This module filters reads based on the minimum readlength and read quality you specify.

**SFilter Reports v1**

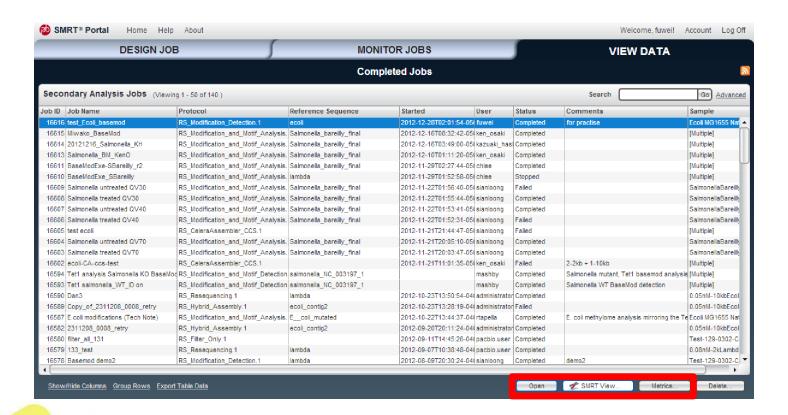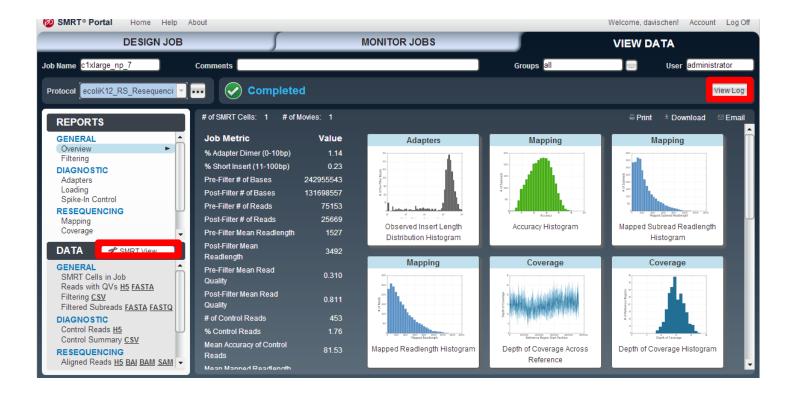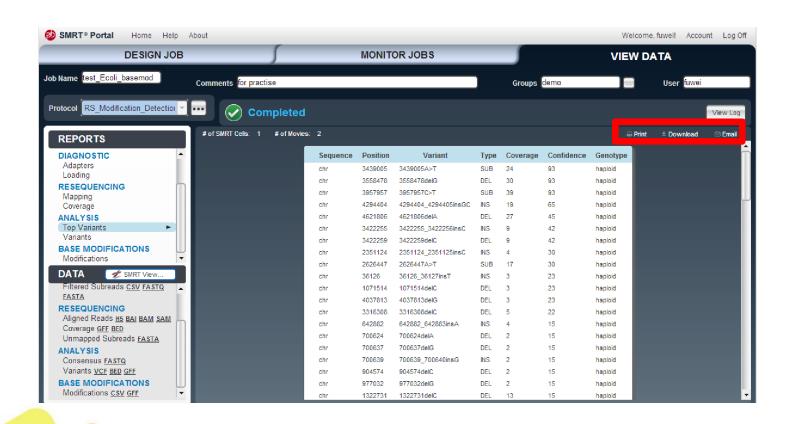This module contains no options.

OK · Apply · Cancel

## LOG

Close Log

[DEBUG] 2012-02-22 10:54:23,143 [SmrtPipeContext 135] Input data is [
file//opt/smrtanalysis/common/inputs_dropbox_automated/2012_02/17/2311160/0024 ]

[INFO] 2012-02-22 10:54:23,143 [SmrtPipeContext 139] Configuration override for PROGRESS_URL: Old: --> New:
https://demo.smrtportal.com:8443/smrtportal/api

[INFO] 2012-02-22 10:54:23,155 [SmrtPipeContext 150] Changing working directory to /mnt/scratch/tmpVHWTKx

[INFO] 2012-02-22 10:54:23,164 [ClusterJMS 124] cluster engine: SGE

template root: /opt/smrtanalysis/analysis/etc/cluster

temporary dir: /shared/smrtanalysis_shared/tmpKYEu53

[INFO] 2012-02-22 10:54:23,200 [Heartbeat 32] heartbeat sleepTime set to 60

[INFO] 2012-02-22 10:54:23,201 [SmrtPipeContext 298] Running uname -a

[DEBUG] 2012-02-22 10:54:23,213 [SmrtPipeContext 303] Output = ['Linux ip-10-140-2-184 2.6.18-238.12.1.el5xen #1
SMP Tue May 31 14:02:29 EDT 2011 x86_64 x86_64 x86_64 GNU/Linux']

[DEBUG] 2012-02-22 10:54:23,213 [SmrtPipeContext 304] Error Code = 0

[DEBUG] 2012-02-22 10:54:23,213 [SmrtPipeContext 305] Error Message =

[INFO] 2012-02-22 10:54:23,213 [SmrtPipeMain 338] smrtpipe running on Linux ip-10-140-2-184 2.6.18-238.12.1.el5xen
#1 SMP Tue May 31 14:02:29 EDT 2011 x86_64 x86_64 x86_64 GNU/Linux

[INFO] 2012-02-22 10:54:23,214 [SmrtPipeMain 351] SMRT Analysis v1.3.0 / SMRTpipe v1.3.0.103818

[INFO] 2012-02-22 10:54:23,214 [SmrtPipeMain 353] Starting smrtpipe v0.9.103324

[INFO] 2012-02-22 10:54:23,354 [HttpProgress 101] Job Progress 'Started' event POSTED to
https://demo.smrtportal.com:8443/smrtportal/api/jobs/016443/status

[DEBUG] 2012-02-22 10:54:23,354 [HttpProgress 102] The data string is

# SMRT Portal

## Kinetic Detections

### Modification Qv Vs. Coverage



### Modification QV Histogram



### Modification QV Histogram

| # of SMRT Cells: 5 | # of Movies: 10 | | | | | | 🖨 Print    ⬇ Download    ✉ Email |
|---|---|---|---|---|---|---|---|
| **Motif** | **Modified Position** | **% Motifs Detected** | **# Of Motifs Detected** | **# Of Motifs In Genome** | **Mean Modification QV** | **Mean Motif Coverage** | **Partner Motif** |
| GATC | 2 | 99.19 | 37929 | 38240 | 132.4 | 76.7 | GATC |
| GCACNNNNNNGTT | 3 | 98.66 | 587 | 595 | 124.0 | 76.6 | AACNNNNNNGTGC |
| AACNNNNNNGTGC | 2 | 98.49 | 586 | 595 | 124.5 | 74.0 | GCACNNNNNNGTT |
| ANCCTGGTCNNK | 3 | 58.33 | 49 | 84 | 59.9 | 77.2 | |
| CCTGGTNNAT | 1 | 40.70 | 105 | 258 | 62.5 | 80.5 | |
| CCTGGYA | 1 | 27.16 | 468 | 1723 | 58.0 | 82.9 | |

# Results Directory



data
log
movie_metadata
results
workflow
index.html
input.fofn
input.xml
job.sh
metadata.rdf
settings.xml
toc.xml
vis.jnlp

controlReport.html
controlReport.xml
coverage.html
coverage.xml
coverageHistogram.png
coverageHistogram_thmb.png
coveragePlot_ref000001.png
coveragePlot_ref000001_thmb.png
filterReports_adapters.html
filterReports_adapters.xml
filterReports_filterStats.html
filterReports_filterStats.xml
filterReports_loading.html
filterReports_loading.xml
mappedRLHistogram.png
mappedRLHistogram_thmb.png
mappedSubreadRLHistogram.png
mappedSubreadRLHistogram_thmb.png

PACIFIC
BIOSCIENCES

http://www.pacb.com/support/software-downloads/

**SMRT Analysis v2.3.0** (released 10/15/2014)

Download

checksum

**SMRT Analysis v2.3.0 Patch 5**

Download

checksum

**Release documentation**

SMRT Analysis Release Notes

SMRT Analysis Software Installation

SMRT Analysis System Requirements

SMRT Portal Help

SMRT View Help

Running SMRT Analysis on Amazon

**Developer documentation**

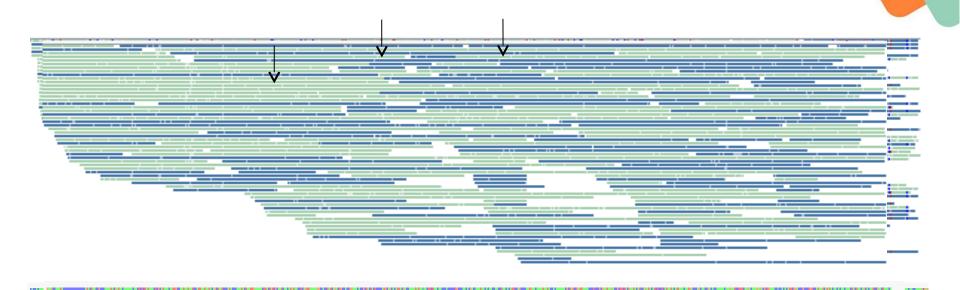SMRT Pipe Reference Guide

Secondary Analysis Web Services API

smrtanalysis 实例

# Hierarchical Genome Assembly Process (HGAp)



1．选定长度作为骨架
2．比对其他 reads
3．矫正错误
4．构建准确的 consensus

- Utilizes every bit of data:
    - Long reads for continuity
    - Shorter reads for improving accuracy
- Accuracy: 85.7% ⛓ 99.3%, 9089 bp
- Chimera / low quality regions can be filtered out early
- Accurate long consensus reads easier to assemble

# SMRTPIPE

# Falcon 实例

DBG2OLC 实例

从源代码编译，使用下面的命令：
g++ -O3 -o SparseAssebmler *.cpp
g++ -O3 -o DBG2OLC *.cpp
g++ -O3 -o Sparc *.cpp

## 选择部分数据

./SelectLongestReads sum 600000000 longest 0 o Illumina_50x.fastq f Illumina_500bp_2x300_R1.fastq
./SelectLongestReads sum 260000000 longest 0 o Pacbio_20x.fasta f Pacbio.fasta
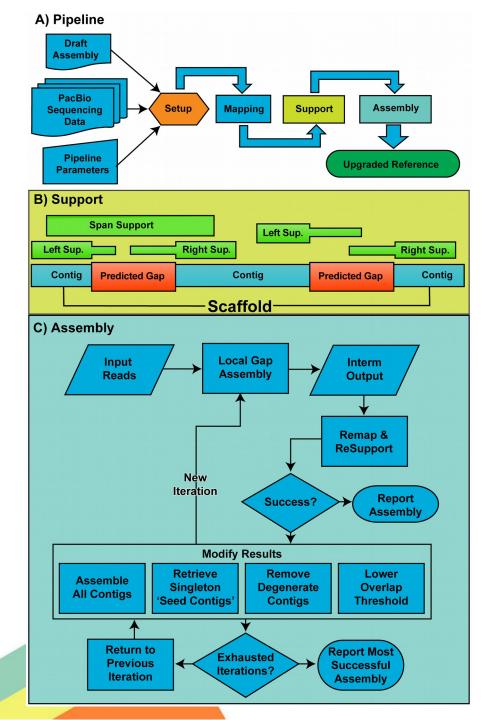
smrtdenovo 实例

1. awk 'NR%4==1||NR%4==2' selfSampleData/pacbio_filtered.fastq | sed 's/^@/>/g' > reads.fa
2. smartdenovo/smartdenovo.pl reads.fa > wtasm.mak
3. make -f wtasm.mak

# PBJelly 实例

**A) Pipeline**

**B) Support**

**C) Assembly**

对于 de novo 组装，长度超过读长的重复序列会产生缺口，导致片段化的组装。很难检测重复区域的变异，而这些对了解某些疾病可能很重要。

PBjelly 能够将 PacBio 长读长序列与组装的草图比对，进行 gapfilling 。研究人员将这种方法应用在四个基因组上，解决了 63%-99% 的 gap

数据： /home/tbc/jellyExample
程序： source /opt/PBSuite_14.1.15/setup.sh

## 1) Create your Protocol.xml

## 2) Run each stage

### Jelly.py &lt;stage&gt; yourProtocol.xml

The stages, in order, and their descriptions are
1. setup        Tag sequence names, find gaps, and index the reference
2. mapping      Use blasr to map the sequences to the reference
3. support      Identify which reads support which gaps
4. extraction   For each gap, consolidate all reads supporting it into a local-assembly folder.
5. assembly     Build the consensus gap-filling sequence
6. output       Stitch the reference sequences and gap-fillling sequences together.

## 3) Passing Parameters through Jelly.py

If you would like to pass a parameter to the stage you are running, use
"-x". For example, when running the support stage, if you only wanted
Jelly to attempt to fill captured-gaps (i.e. no inter-scaffold gaps), and
you wanted to require that a read must have a minimum mapping QV of >=
250 to support a gap, you'd use the command:
> Jelly.py support Protocol.xml -x "--capturedOnly --minMapqv=250"
All parameters you add need to be enclosed in double quotes after the -x

Isoseq 分析

命令行运行 ReadsOfInser 的命令是 ConsensusTools．可以到 smrt 安装目录 <smrtanalysis_directory>/doc/bioinformatics-tools/ConsensusTools/doc/index.html 下面查看相关文档。

需要的输入文件：

```
input.fofn --- A plain file containing the list of .bax.h5 file locations.
```

下面是一个 input.fofn 文件的例子 . 一共有两个
movie ，每个包含三个 .bax.h5 文件

```
/MYHOME/runs/m140121_100730_42141_c100626750070000001823119808061462_s1_p0.1.bax.h5
/MYHOME/runs/m140121_100730_42141_c100626750070000001823119808061462_s1_p0.2.bax.h5
/MYHOME/runs/m140121_100730_42141_c100626750070000001823119808061462_s1_p0.3.bax.h5
/MYHOME/runs/m140121_132657_42141_c100626060070000001823118408061490_s1_p0.1.bax.h5
/MYHOME/runs/m140121_132657_42141_c100626060070000001823118408061490_s1_p0.2.bax.h5
/MYHOME/runs/m140121_132657_42141_c100626060070000001823118408061490_s1_p0.3.bax.h5
```

## 示例命令：

```
ConsensusTools.sh CircularConsensus  --minFullPasses 0  --minPredictedAccuracy 75 \
   --parameters <smrtanalysis_directory>/analysis/etc/algorithm_parameters/2014-03 \
   --numThreads 24 --fofn /MYHOME/test_dir/input.fofn \
   -o /MYHOME/test_dir/data
```