#### NGS解析をはじめよう

~よりよいデータ解析のためのFASTQファイルの前処理と



For Research Use Only. Not for use in diagnostic procedures.

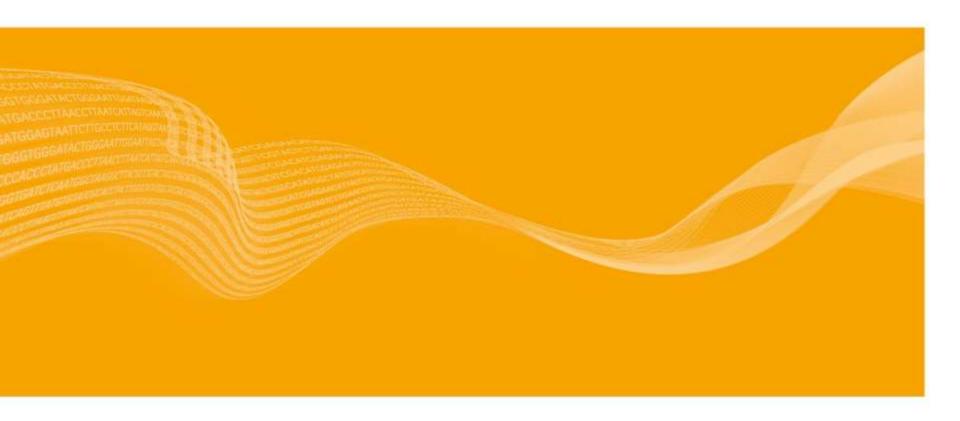


#### ウェビナーの概要

- データ解析ワークフロー FASTQファイル作成までとその後の2次解析)
  - データ解析フローの概要
  - FASTQ ファイルフォーマット
- FASTQ ファイル前処理の種類
  - アダプタートリミング
  - クオリティトリミング
  - リードの結合
  - サブサンプリング
- FASTQ ファイル前処理の具体例
  - Fastqファイル作成と同時にアダプタートリミングする方法
  - BaseSpace Sequence Hub (BSSH) FASTQ Toolkit App を用いた前処理
- FASTQ ファイルのクオリティ評価
  - BaseSpace Sequence Hub (BSSH) FastQC App を用いたクオリティ評価

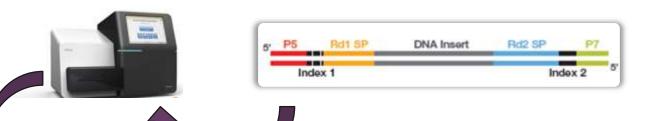


## データ解析のワークフロー





#### データ解析フロー



- ・シーケンス 画像取得 / シグナル抽出
- ベースコール (塩基の決定)

\*.bcl ファイル

#### FASTQ ファイル

- 装置内ソフトウェア (Local Run Manager、 MiSeq Reporter)
- BaseSpace Sequence Hub (BSSH)
- Bcl2fastq (Linux)

- 2次解析ワークフロー
- Third Party

BAMs, VCFs, Assembly



#### FASTQ ファイル

ファイル名
SampleName\_SampleNumber\_Lane\_Read\_FlowCellIndex.fastq.gz
SampleName\_S1\_L001\_R1\_001.fastq.gz
SampleName\_S1\_L001\_R2\_001.fastq.gz

FASTQ ファイルフォーマット

@SIM:115:FCX:4:2106:6329:1045 ATGCTAGCC 1:N:0:ATCCGA へッダー
TCGCACTCAACGCCCTGCATATGACAAGACAGAATC 塩基配列
+
<>;BB=><9=AAAAAAAAAAAAB8:<B<;<<<????B= クオリティスコア

ヘッダー名

@Instrument:RunID:FlowCellID:Lane:Tile:X:Y:UMI Read:Filter:0:IndexSequence or SampleNumber

2020/03/25

NGS解析をはじめよう ~基礎からわかるバイオインフォマティクス入門編~テクニカルアプリケーションサイエンティスト 仁田原 翔太 https://jp.illumina.com/events/webinar/2020/webinar-0325-j.html



#### クオリティスコア (Qスコア)

- クオリティスコアは、各塩基のベースコールの確からしさを表す指標。データ評価を行う時に用いる。
  - Q10:90%の確率でそのベースコールは確からしい
  - Q20:99%の確率でそのベースコールは確からしい
  - Q30:99.9%の確率でそのベースコールは確からしい
  - Q40: 99.99 %の確率でそのベースコールは確からしい

Q-Score	Symbol	Q-Score	Symbol	Q-Score	Symbol	
0	!	14	/	28	=	
1	í.	15	0	29	>	
2	#	16	1	30	?	
3	\$	17	2	31	@	
4	%	18	3	32	А	
5	&	19	4	33	В	
6	•	20	5	34	С	
7	(	21	6	35	D	
8	)	22	7	36	E	
9	*	23	8	37	F	
10	+	24	9	38	G	
11	,	25	:	39	Н	
12	-	26	;	40		
13		27	<			



### データ解析のワークフロー ~FASTQファイルが手に入ったら~

#### FASTQ ファイル

- 装置内ソフトウェア (Local Run Manager、MiSeq Reporter)
- BaseSpace Sequence Hub (BSSH)
- Bcl2fastq (Linux)

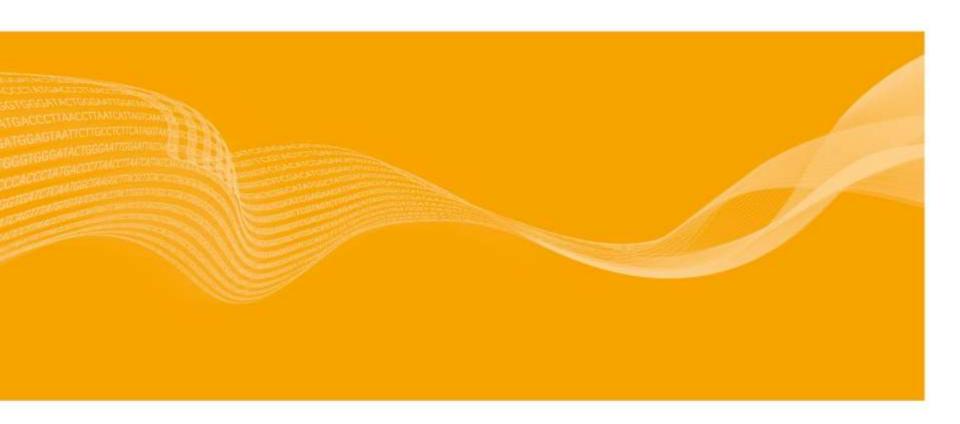
- ① FASTQファイルの特徴 を知る!(データ評価)
- → クオリティスコア分布 リード長分布
- → FastQC App
- ② FASTQファイルの 前処理
  - → FASTQ Toolkit App
- ③ FASTQファイルの確認
  - → FastQC App

# BAMs, VCFs, Assembly

- 2次解析ワークフロー
- Third Party

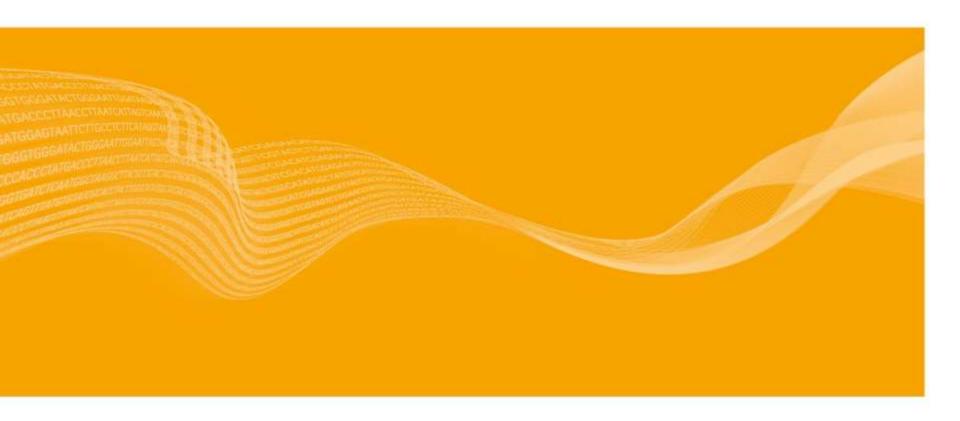


#### FASTQファイル前処理の種類





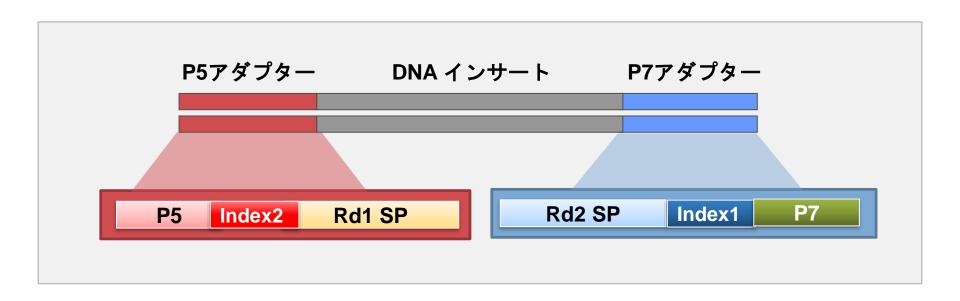
### アダプタートリミング





## ライブラリ構造 ーアダプターとは?ー

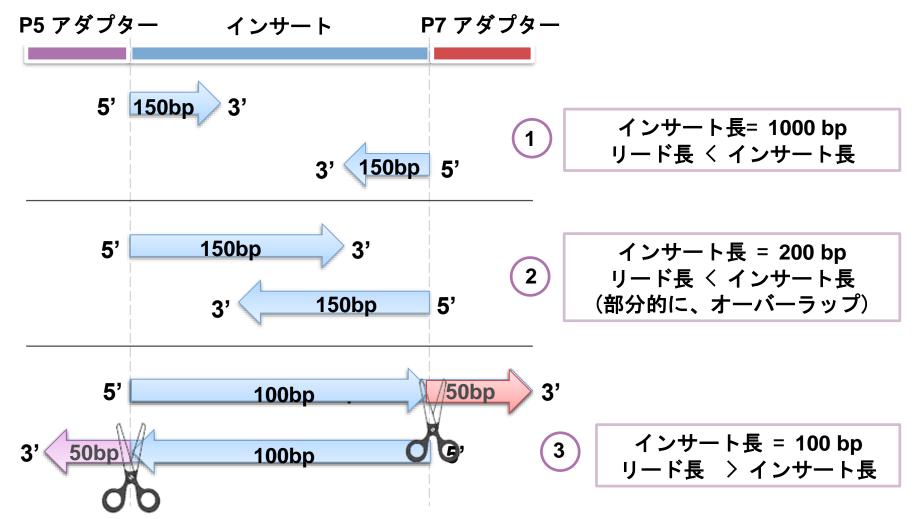
- ライブラリ= DNA インサート + アダプター
- Read 1 シーケンスプライマー (Rd1 SP)
- Read 2 シーケンスプライマー (Rd2 SP)





#### アダプタートリミング

例:150 bp x 2



# アダプタートリミング(ハードトリミング)



@M00000:71:000000000-D00LW:1:1101:16265:1658 1:N:0:1

ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCCCTCTCCCTTATACACATCTCCGAGCCCA

+

@M00000:71:00000000-D00LW:1:1101:16265:1658 1:N:0:1

ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCC

+

BCCCCFFCCBCCGGGGGGGGGGGGGHHHHHHHHHHHHHH



アダプタ一配列以降を完全に除去する



# アダプタートリミング (アダプターマスキング)

@M00000:71:00000000-D00LW:1:1101:16265:1658 1:N:0:1

ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCCCTCTCTTATACACATCTCCGAGCCCA

+



@M00000:72:00000000-D00LW:1:1101:16265:1658 1:N:0:1

+

アダプター配列をNでマスキングし、クオリティスコアを2 (#) に置き換える (もし、後続のデータ解析で全て同じリード長の配列が必要な時に使用)



### どうしてアダプタートリミングが必要?



#### Higher alignment %



# BWA (backtrace)

Sample	Sample Name	Total Aligned Reads	Percent Aligned Reads
1	NA12892	354,882	77.4%
2	NA12892- trim	450,007	98.2%



#### どうしてアダプタートリミングが必要?





#### Improved assemblies

Data: 2 x 250bp, *E.coli* ゲノムサイズ: 4.64M (Nextera™ XT)

Assembly metrics		Before adapter trimming	After adapter trimming	
	N50 21		29,791	
Maximum contig		553	174,326	
Assembly length		18,497,207	4,876,437	
Number of contigs		1,387,508	1,115	

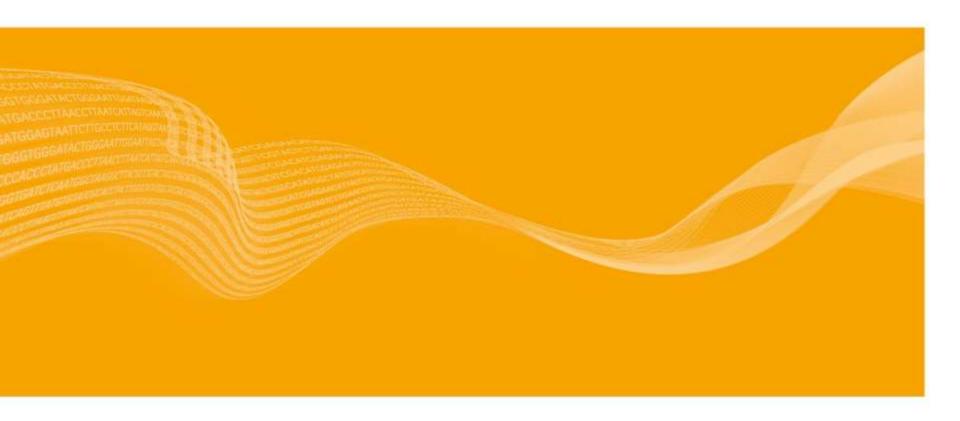
N50: アッセンブリした配列を長い順に並べて、長い方から順に配列の塩基数を

足していった時に全体の長さの半分に達したときの配列の塩基数(bp)

Number of contig: アッセンブリしてつなぎ合わされた配列の数



## クオリティトリミング





#### クオリティトリミング

- 平均クオリティスコアを基に、リード両末端のクオリティの低い 塩基を取り除く
- どんな時にトリミングする?
  - 基本的に、全てのアプリケーション
  - 特に3'末端のクオリティの低い塩基が、後続の解析にクリティカルに 影響を与える時 (e.g. de novo アセンブリ, リード結合など)
- どんな時にトリミングしない?
  - リシーケンス解析: ほとんどのソフトウェア (i.e. BWA, Isaac) が クオリティスコアを考慮に入れてアライメントしているため。
- イルミナツール
  - BSSH FASTQ Toolkit App



#### クオリティトリミングの例

#### QualityScoreTrim,20

@M00000:72:00000000-D00LW:1:1101:22420:18334 1:N:0:1

CACCAAGGGCCTGGGGTGTCAATGGCGGGGCTTGTGACTGCACAAAAGGGGCCTCCCGCAGGGGCTCCCGCC

+



@M00000:72:00000000-D00LW:1:1101:22420:18334 1:N|:0:1

CACCAAGGGCCTGGGGTGTCAATGGCGGGGCTTGTGACTGCACAAAAGG

+

BBBBBBBBBBBGGGGEEFGGGHHHHGGG00>10B355@BB3@3BG1?E



Q	ASC
13	
14	/
15	0
16	1
18	3
20	5
22	7
24	9
30	?
31	@
32	Α
33	В
36	Е



#### クオリティトリミングの例

#### QualityScore:20, WindowLength:4

 $20 \times 4 = 80 > 75$ 

@M00000:72:000000000-D00LW:1:1101:22420:18334 1:N:0:1

CACCAAGGGCCTGGGGTGTCAATGGCGGGGCTTGTGACTGCACAAAAGGGGCCTCCCGCAGGGGCTCCCGCC

@M00000:72:000000000-D00LW:1:1101:22420:18334 1:N:0:1

CACCAAGGGCCTGGGGTGTCAATGGCGGGGC

+

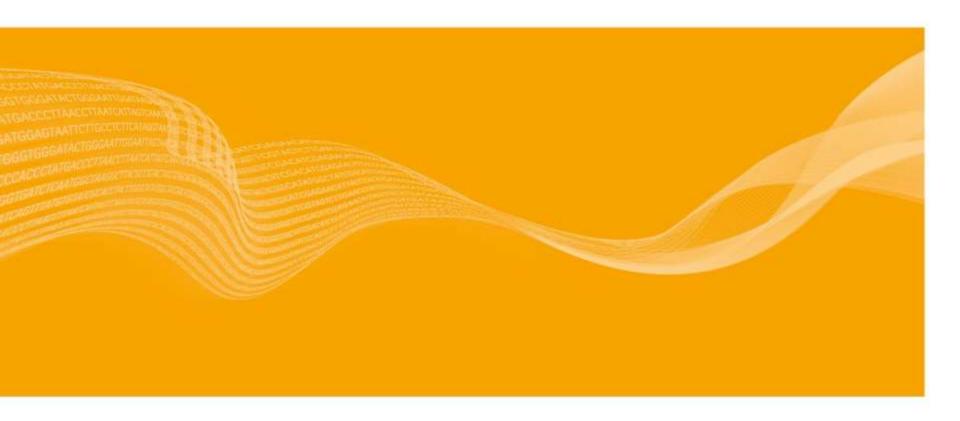
BBBBBBBBBBBGGGGEEFGGGHHHHGGG00



Q	ASC
13	
14	/
15	0
16	1
18	3
20	5
22	7
24	9
30	?
31	@
32	А
33	В



# リード結合



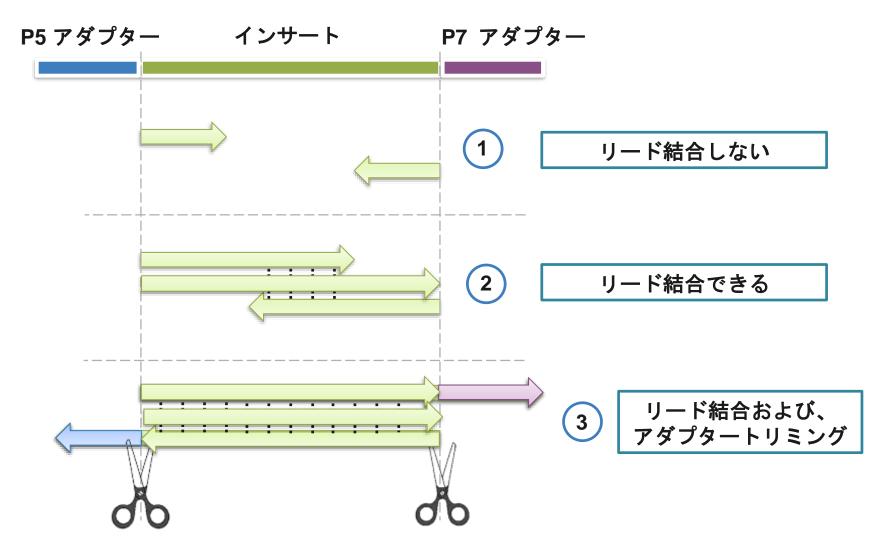


#### リード結合

- ペアードエンドのR1とR2をつなぎ合わせ、シングルリードに すること
- どんな時にリード結合が必要?
  - より長いリードを得ることが重要な時
  - シングルリードしか受け付けないソフトウェアを使用する時
  - 大部分のリードがオーバーラップしている時
  - オーバーラップ領域にindelが検出される時
- どんな時にリード結合が必要ではない?
  - 多くのリードにオーバーラップ領域がないとき
  - オーバーラップ領域に繰り返し配列があるとき
  - 後続の解析で、R1、R2ファイルを別々にインプットする必要があるとき
- 一般的に用いられるツール
  - USEARCH, FLASH



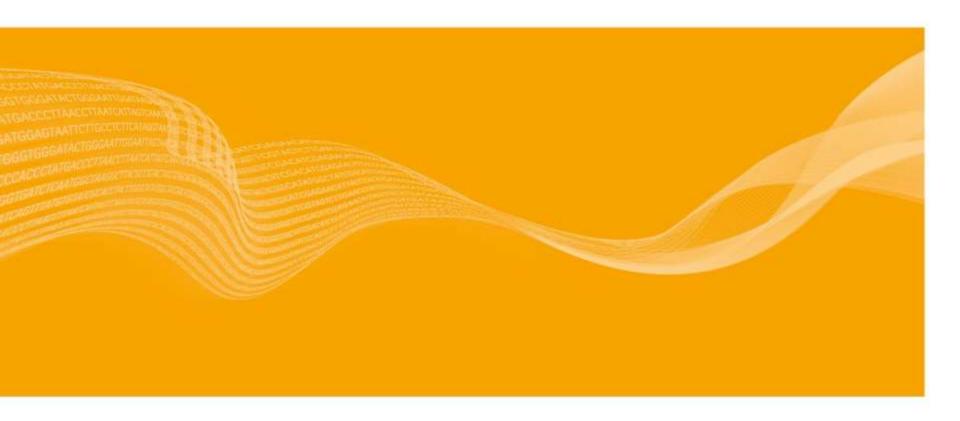
#### リード結合



For Research Use Only. Not for use in diagnostic procedures.



#### サブサンプリング



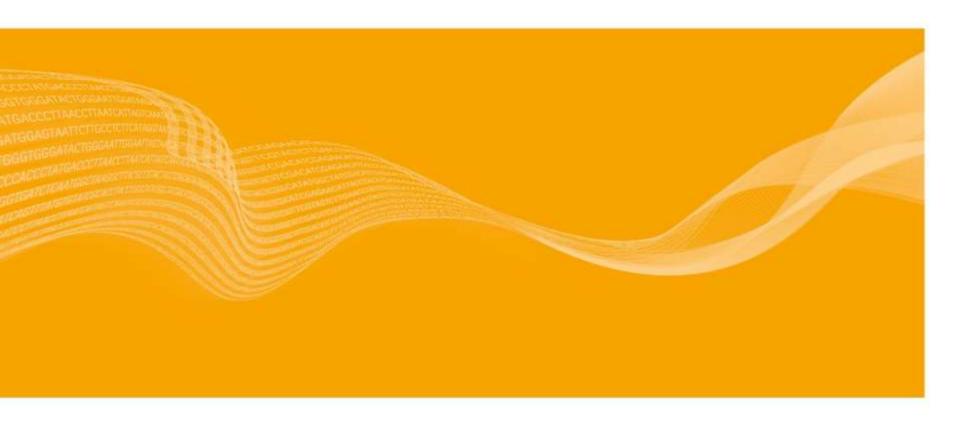


#### サブサンプリング

- シーケンスリード数が多すぎるときに、ランダムに一部の リードのみを抜き出してサブサンプルセットを作成すること
- どうしてサブサンプルが必要?
  - ▶ トラブルシューティングのための解析をより短い時間で行うため
  - パソコンの性能が十分ではなく、シーケンスデータが大きすぎて、 データ解析が途中でタイムアウトする場合
  - イルミナBaseSpace Sequence Hub のいくつかのアプリには、 インプット量に制限があるため
  - 解析に適当なリード数の検討を行う場合
- イルミナツール
  - BSSH FASTQ Toolkit App

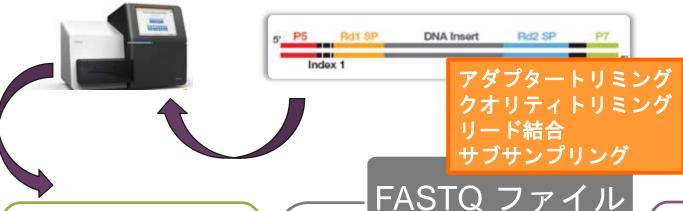


### FASTQファイルの前処理の方法 一具体例一





#### データ解析フロー



- ・シーケンス 画像取得 / シグナル抽出
- ベースコール (塩基の決定)

\*.bcl ファイル

FASIQ J711

- 装置内ソフトウェア (Local Run Manager、 MiSeq Reporter)
- BaseSpace Sequence Hub (BSSH)
- Bcl2fastq (Linux)

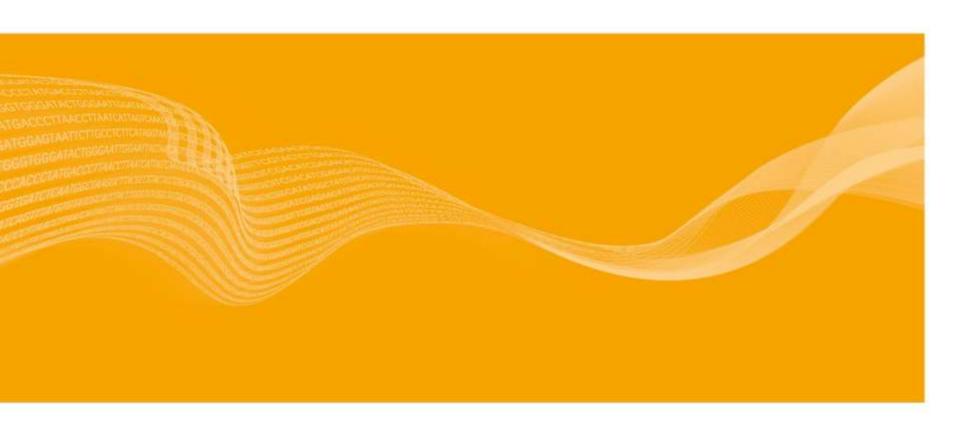
- 2次解析 ワークフロー
- Third Party

BAMs, VCFs, Assembly

アダプタートリミング



#### FASTQファイル作成時にアダプタートリミングする方法





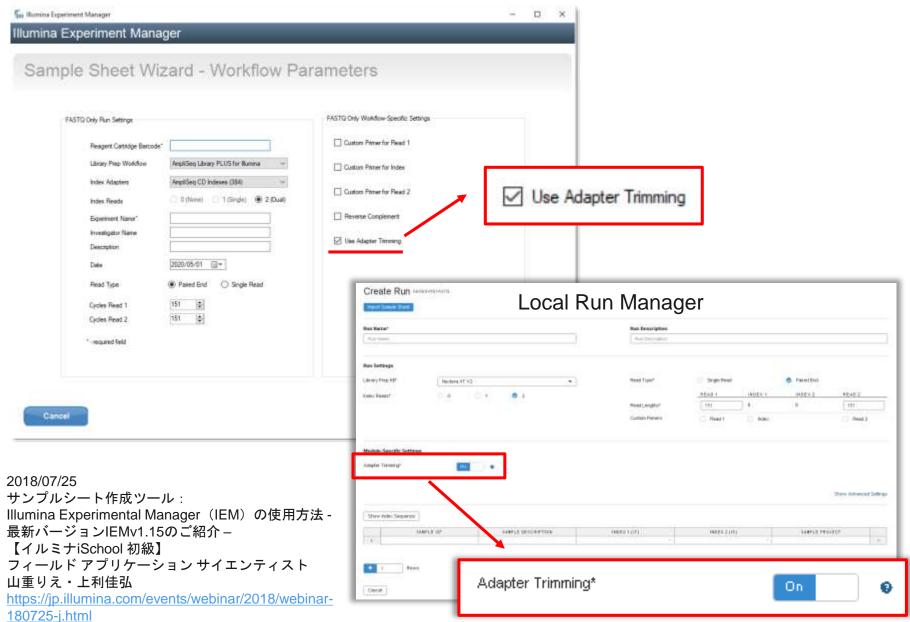
#### アダプタートリミング

# MiSeq™ Reporter, Local Run Manager, BaseSpace™ Fastq Generation and Bcl2fastq

[Header] IEMFileVe 4 Date 4/11/2017	Hard trimr (シーケンス	ming からアダプタ	一配列で	を完全に	取り除	<)
Workflow GenerateFASTQ Applicatic FASTQ Only Assay Nextera XT Description	Read1とRead2のシーケンスが異なる場合は、 以下のようにそれぞれの配列を入力する					
Chemistry Amplicon  [Reads]  151  151	[settings] Adapter, AdapterRea	 ad2,				
[Settings] Adapter CTGTCTCTTATACACATCT						
[Data] Sample_ICSample_Nam Sam Test		7_Index_ index I701 TAAGGCG.	I5_Index_ S502	index2 CTCTCTAT	· -	Description

https://support.Illumina.com/bulletins/2016/12/what-sequences-do-i-use-for-adapter-trimming.html For Research Use Only. Not for use in diagnostic procedures.

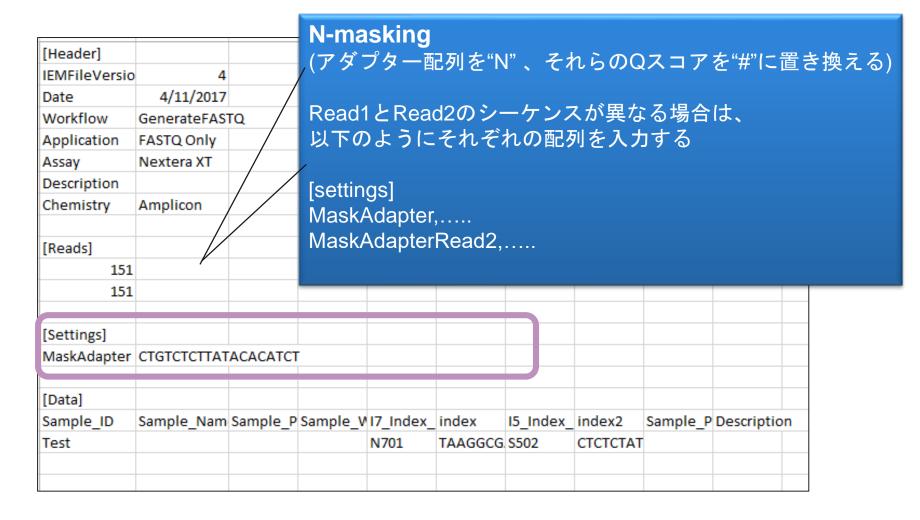






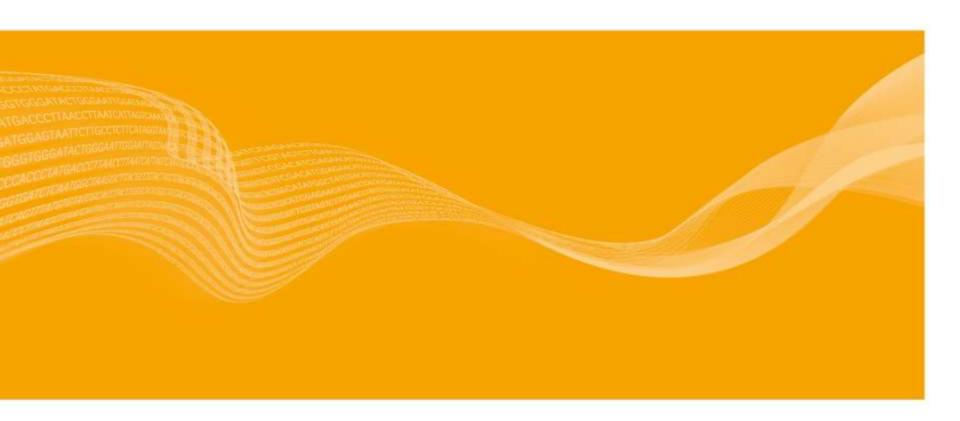
# アダプターマスキング

#### Bcl2fastq





# BaseSpace Sequence Hub FASTQ Toolkit Appを用いた FASTQファイルの前処理









#### **BaseSpace Sequence Hub**

- BaseSpace Sequence Hub (BSSH) は、イルミナ社が提供するクラウドサービスで、イルミナ次世代シーケンサーのデータ管理及び解析にご利用いただけます(有償)。
- 90種類以上のアプリをご提供。解析結果は「project」に保存。
- アプリには、Core App、 BaseSpace Lab App、Third Party Appの3種類。

#### **Core App**



#### BaseSpace Lab App





**Third party App** 



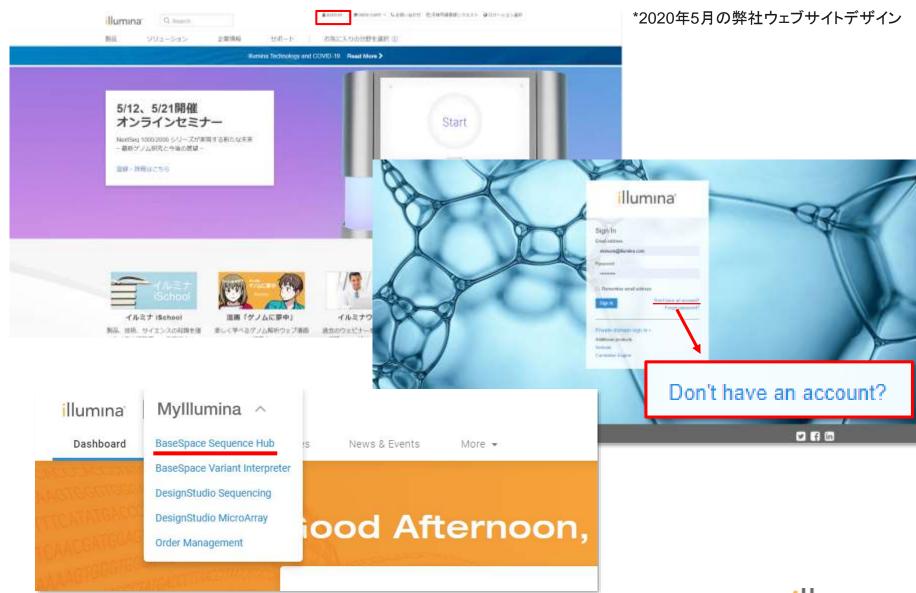


BaseSpace Sequence Hubを使用したデータ解析の基本【イルミナiSchool 初級】 テクニカル アプリケーション サイエンティスト 渡邊 大





#### BaseSpace Sequence Hubのアカウント取得

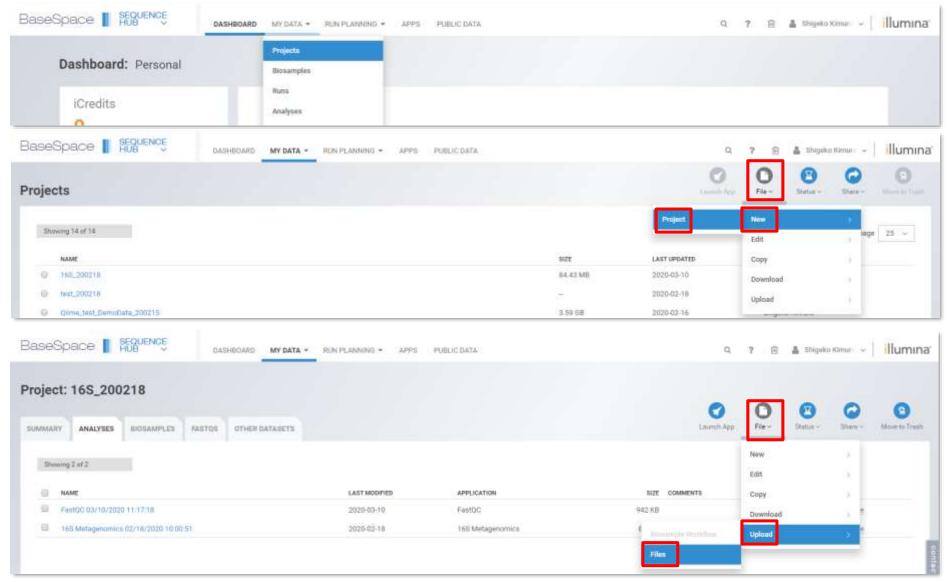


#### BaseSpace Sequence Hub App (>90種類)



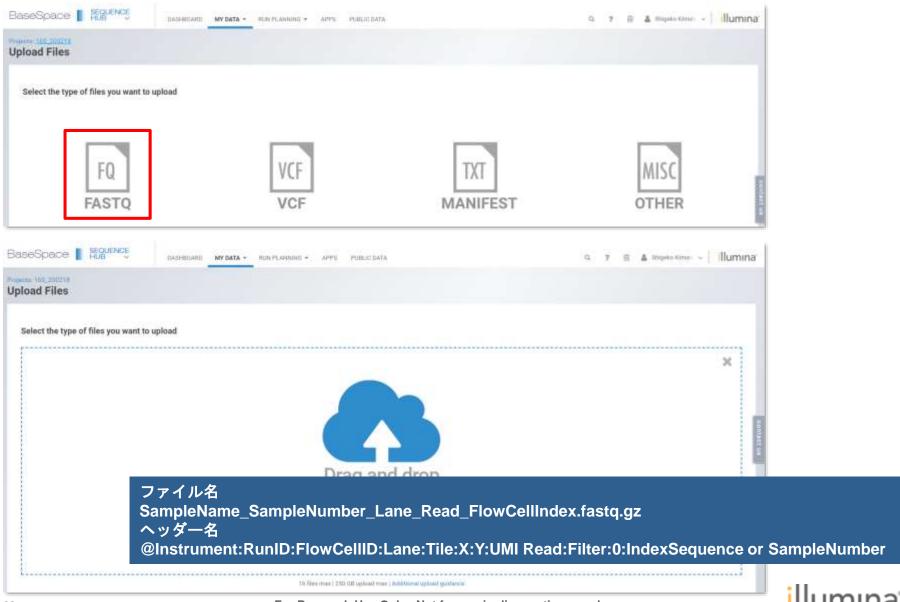


#### FASTQファイルのアップロード法



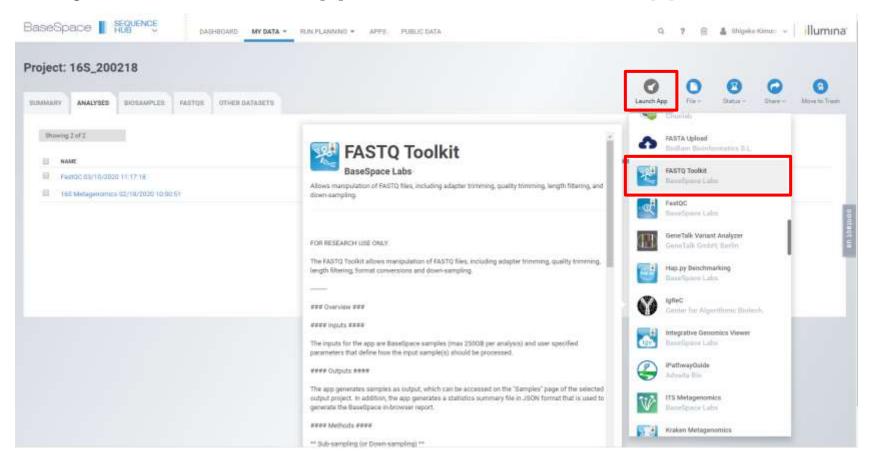


#### FASTQファイルのアップロード法



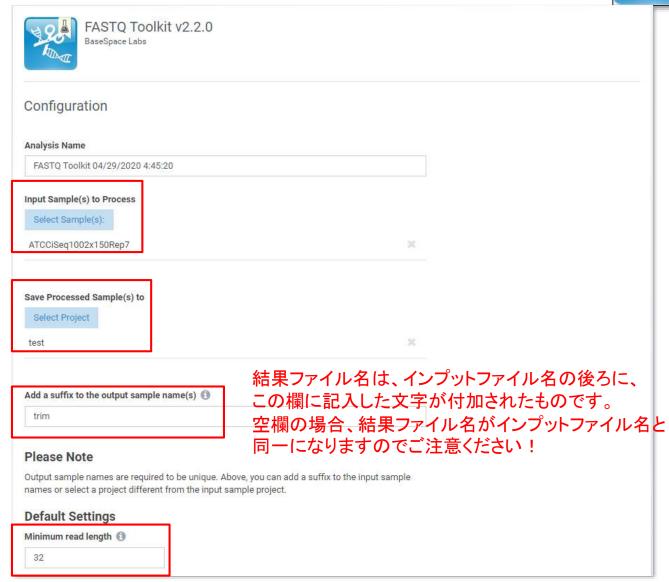


- App→ "FASTQ Toolkit" Appアイコンを選択
- Project内、Launch app→"FASTQ Toolkit" Appを選択











#### ● 適当なOptionを選択

#### Sub-sampling

Sub-sampling is required when only a subset of the sample can be processed by an application (e.g. de novo assembly with memory constrains) or it is not necessary to process a full sample (e.g. for validating an approach at varying levels of genomic coverage). Samples can be sub-sampled to a user defined number of reads or number of bases, or to a specified fraction of the input sample (e.g. 50% of input sample) - again, based on reads or bases.

+ Click to

Click to expand/hide sub-sampling settings

#### Adapter Trimming

Use these parameters if you want to trim adapter sequences. The adapter sequence can be specified separately for the 5'- and 3'-end and is a required input for adapter trimming. In addition to the adapter sequence, you can specify a stringency value that determines how similar a sequence has to be in order to be identified as an adapter sequence.

+ Click to expand/hide adapter trimming settings

#### **Base Trimming**

Bases can be trimmed from either the 5'- or 3'-end by specifying a number of bases for each end. Alternatively, bases can be trimmed from the 3'-end by specifying a maximum read length after trimming (e.g. to trim 150bp reads to 100bp reads).

+ Click to expand/hide base trimming settings

#### **Quality Trimming**

Use these parameters to trim low quality bases. Please note that most aligners (including BWA and Isaac) perform quality trimming internally during alignment. Quality trimming is more appropriate for workflows such as assembly, metagenomics or Methyl-Seq.

+ Click to expand/hide quality trimming settings

#### Poly-A/T Trimming

Poly-A/T tails are considered repeats of As or Ts at the read ends. The minimum length of a tail can be specified below. A small number of tails can occur even after trimming poly-A/T tails. For example, a sequence that ends with AAAAATTTTT and that has been trimmed for the poly-T will still contain the poly-A. Trimming poly-A/T tails can reduce the number of false positives during database searches, as long tails tend to align well to sequences with low complexity or sequences with tails (e.g. viral sequences) in the database. It also improves assemblies of (meta-)transcriptomes.







Adapter Trimming optionを選択

lapter Trimming	None selected		
e these parameters if you want to trim adapt	Nextera Ranid Canture (CTGTCTCTTATACACATCT)		
ecified separately for the 5'- and 3'-end and is adapter sequence, you can specify a string			
s to be in order to be identified as an adapte			
<ul> <li>Click to expand/hide adapter tri</li> </ul>	Small RNA v1 (TCGTATGCCGTCTTCTGCTTGT)		
Click to expand/mide adapter un	Small RNA v1.5 (ATCTCGTATGCCGTCTTCTGCTTG)		
Adapter trim stringency (0.01-0.99) 📵	TruSeq Small RNA (TGGAATTCTCGGGTGCCAAGG)		
0.90	TruSeq Dual Index (AATGATACGGCGACCACCGAGATCTACAC,GATCGGAAGAGCACACGTCTGAACTCCAGTCAC		
	TruSeq HT/LT (AGATCGGAAGAGCACACGTCTGAACTCCAGTCA,AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTC		
Select an adapter to trim or specify the	TruSeq HT/LT Common Sequence (AGATCGGAAGAGC)		
None selected	▼		
1-			
Adapter sequence(s) to trim from the 5'	-end 📵		
**	***		
Adapter sequence(s) to trim from the 3'	-end 🕦		
Trim Ns from the 3'-end before identifyi	ng adapters 🕦		
	:11:		





● Quality Trimming optionを選択

#### **Quality Trimming**

Use these parameters to trim low quality bases. Please note that most aligners (including BWA and Isaac) perform quality trimming internally during alignment. Quality trimming is more appropriate for workflows such as assembly, metagenomics or Methyl-Seq.

Click to expand/hide qualit	ty trimming settings	
Trim bases at the 5'-end with a q	The state of the s	+11= , 48
	5'末端の読み始めでク 悪い塩基が続く場合	イオリティか
Trim bases at the 3'-end with a q	51. 51.	5カナリニノボ
	3'末端のリード後半で 悪い塩基が続く場合	ジオリティか
Trim the 3'-end of reads with que	ality score ①	٦
20		
Trim 3'-end using a sliding windo	ow approach with window length 📵	Sliding Windowの設定
4		





● Sub-Sampling optionを選択

#### Sub-sampling Sub-sampling is required when only a subset of the sample can be processed by an application (e.g. de novo assembly with memory constrains) or it is not necessary to process a full sample (e.g. for validating an approach at varying levels of genomic coverage). Samples can be sub-sampled to a user defined number of reads or number of bases, or to a specified fraction of the input sample (e.g. 50% of input sample) - again, based on reads or bases. Click to expand/hide sub-sampling settings Maximum number of FASTQ entries to keep (1) 絶対値 残したい最大リード数の設定 Maximum percentage of FASTQ entries to keep [] 割合 Maximum number of bases to keep 🕦 絶対値 残したい最大塩基数の設定 Maximum percentage of bases to keep 📵 割合 Choose reads at random when sub-sampling instead of taking the first n reads 📵





条件設定が終わったら、規約に同意し、Launch Applicationを 選択

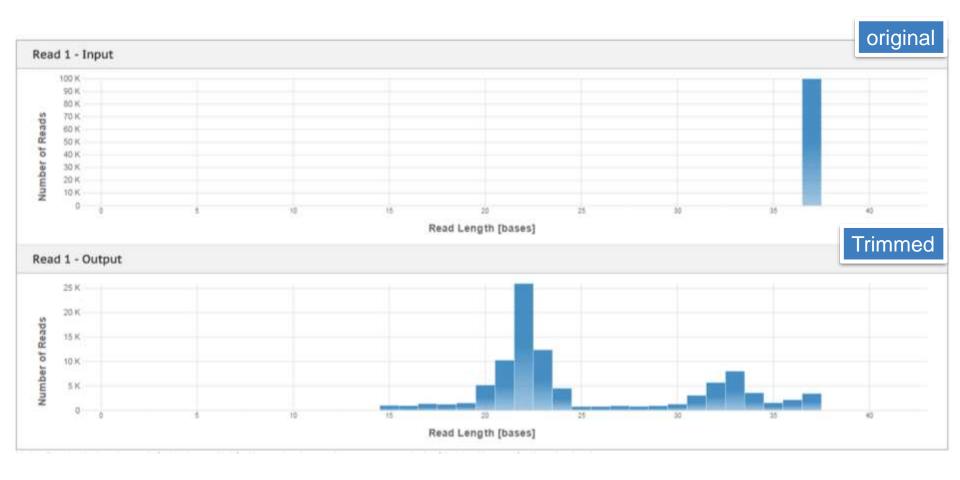
#### BaseSpace Labs Disclaimer

I acknowledge and agree that (i) this is a BaseSpace Labs App, (ii) I am using it AS-IS without any warranty of any kind, (iii) Illumina has no obligation to provide any technical support for this App, and (iv) Illumina has no liability for my use of this App, including without limitation, any loss of data, incorrect results, or any costs, liabilities, or damages that result from use of this App.

Launch Application



# **FASTQ Toolkit output- Results**





# FASTQ Toolkit その他のオプション

リード末端 の Poly-A/T を除去

Poly-A/T Trimming

5'あるいは3'末端から、 指定した長さの塩基を 除去

**Base Trimming** 

ヘッダーライン 等の修正

Fix Format

リード長、平均クオリティ スコア、GC含有率を指定し、 条件を満たさないリードを 除去

Read Filtering

ペアエンドをシングルエンド に変換、 メイトペアリードを

ペアエンドリードに変換

Modify Reads



# データ解析のワークフロー ~FASTQファイルが手に入ったら~

#### FASTQ ファイル

- 装置内ソフトウェア (Local Run Manager、MiSeq Reporter)
- BaseSpace Sequence Hub (BSSH)
- Bcl2fastq (Linux)

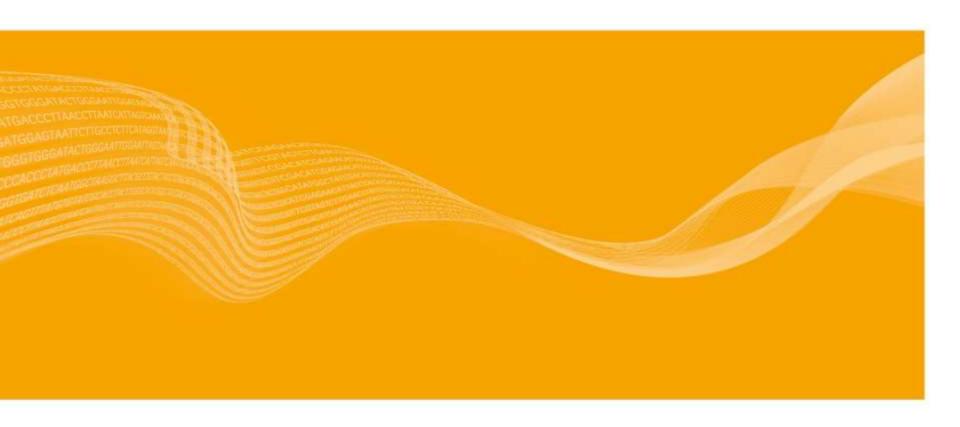
- ① FASTQファイルの特徴 を知る!(データ評価)
- → クオリティスコア分布 リード長分布
- → FastQC App
- ② FASTQファイルの 前処理
  - → FASTQ Toolkit App
- ③ FASTQファイルの確認
- → FastQC App

#### BAMs, VCFs, Assembly

- 2次解析ワークフロー
- Third Party



## シーケンスデータのクオリティ評価





### 私のシーケンスデータのクオリティはいい?

#### **BSSH FastQC App**



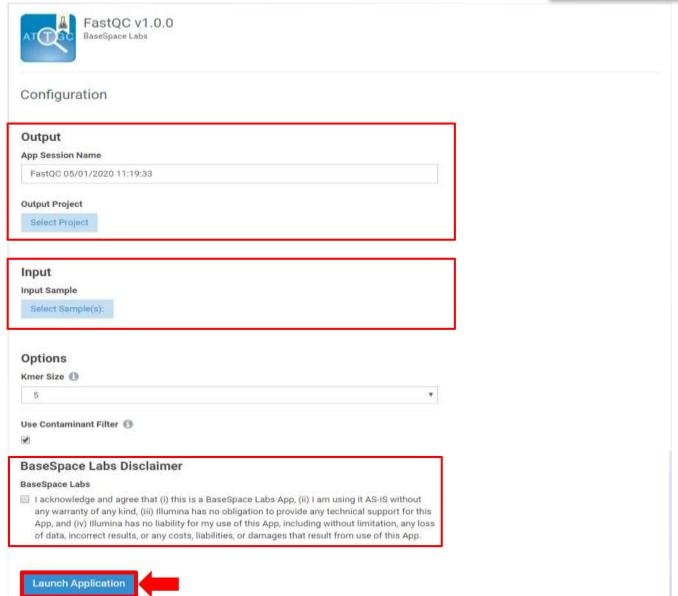
- 配列データのクオリティ評価を行うためのApp (サブサンプリングしたデータを使用すればより早く行える)
- 非常に様々なグラフが結果として出される。アプリケーションによって、チェックしなければならない項目が異なるので、いずれかの結果が"Failure ②"だからと言って、クオリティが悪いということではない。
- ライブラリの特徴を知ることができる
  - クオリティスコアの分布
  - リード長の分布(アダプタートリミングされている?)

過去のデータと比較



## FastQC 入力画面







# FastQC: 基本情報



• 解析の基本情報を表で確認できる

ファイル名、全リード数、リード長、GC含量

В	Basic Statistics				
	Measure	Value			
	Filename	ATCCiSeq1002x150Rep7_S7_L001_R2_001.fastq.gz			
	File type	Conventional base calls			
	Encoding	Sanger / Illumina 1.9			
	Total Sequences	314069			
	Filtered Sequences	0			
	Sequence length	151			
	%GC	52			

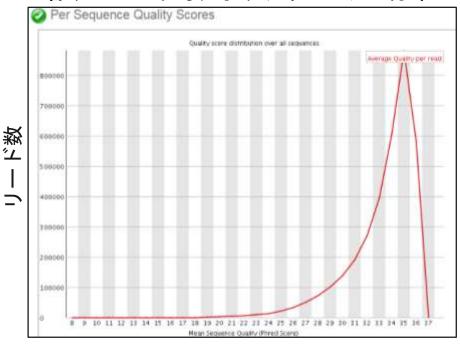


#### **FastQC:**

## クオリティスコアの分布



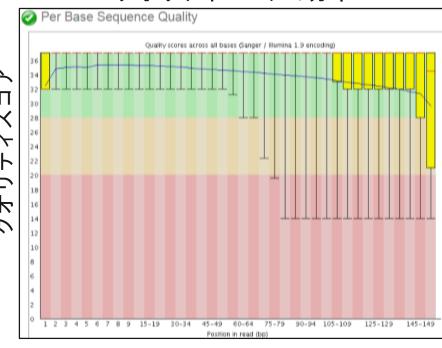
#### 各リードの平均クオリティスコアの分布



リードの平均クオリティスコア

平均クオリティスコアのピークはQ35 ほとんどがQ30以上 →クオリティの良いデータが得られている

クオリティスコアの分布



リード中の塩基の位置(bp)

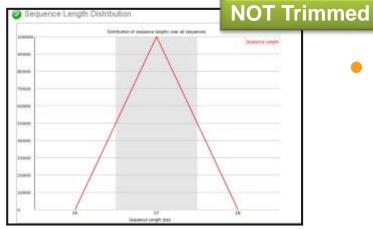
クオリティスコアの分布 3'末端のクオリティスコアが落ちている →クオリティトリミングを要検討



# FastQC: リード長の分布

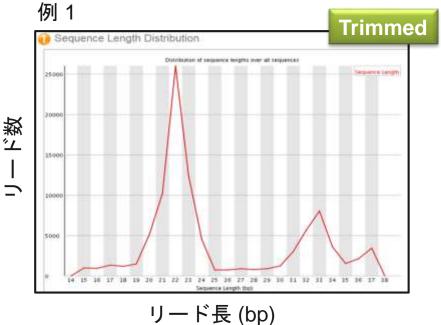


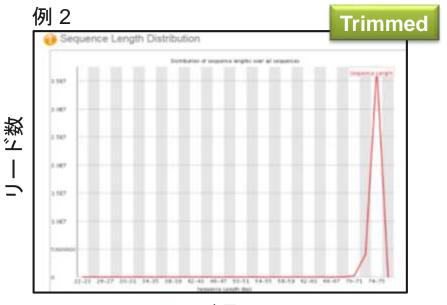
リード数



トリミングされていなければ、 1つの大きなピークが見られる

リード長 (bp)





リード長 (bp)

## 参考資料

- FASTQ files explained
- What sequences do I use for adapter trimming?
- Adapter trimming: Why are adapter sequences trimmed from only the 3' ends of reads?
- FASTQC detailed documentation
- 【ウェビナー】NGS解析をはじめよう ~基礎からわかるバイオイン フォマティクス入門編~
- 【ウェビナー】サンプルシート作成ツール: Illumina Experimental Manager (IEM) の使用方法 - 最新バージョンIEMv1.15のご紹介 -
- 【ウェビナー】 BaseSpace Sequence Hubを使用した データ解析の基本



## ご清聴ありがとうございました

